

Logistic Regression and Softmax Regression

Softmax regression is a generalization of logistic regression for multiclass

Logistic regression is softmax regression when $K=2$

Proof:

Logistic regression states that:

$$P(Y=1|X, \theta) = \sigma(X, \theta_1) = \frac{1}{1 + e^{-(\theta_{01} + \theta_1^T X)}} = \frac{e^{\theta_{01} + \theta_1^T X}}{1 + e^{\theta_{01} + \theta_1^T X}}$$

$$P(Y=0|X, \theta_0) = 1 - P(Y=1|X, \theta_1)$$

$$\text{if } Y=0, 1 = e^{\theta_{00} + \theta_0^T X}$$

$$= 1 - \frac{e^{\theta_{10} + \theta_1^T X}}{e^{\theta_{00} + \theta_0^T X} + e^{\theta_{10} + \theta_1^T X}}$$

$$= \frac{e^{\theta_{00} + \theta_0^T X} + e^{\theta_{10} + \theta_1^T X}}{e^{\theta_{00} + \theta_0^T X} + e^{\theta_{10} + \theta_1^T X}}$$

$$= \frac{e^{\theta_{00} + \theta_0^T X}}{e^{\theta_{00} + \theta_0^T X} + e^{\theta_{10} + \theta_1^T X}}$$

$$= \frac{e^{\theta_{10} + \theta_1^T X}}{e^{\theta_{00} + \theta_0^T X} + e^{\theta_{10} + \theta_1^T X}}$$

roughly:

$$P(Y=1|X, \theta) = \frac{e^{\theta_1^T X}}{e^{\theta_0^T X} + e^{\theta_1^T X}}$$

$$P(Y=0|X, \theta) = \frac{e^{\theta_0^T X}}{e^{\theta_0^T X} + e^{\theta_1^T X}}$$

Softmax regression states that:

for $K=2$:

$$z_1 = \theta_1^T X$$

$$z_2 = \theta_2^T X$$

Softmax is:

$$\frac{1}{\sum_{j=1}^K e^{z_j}} \begin{bmatrix} e^{z_1} \\ \vdots \\ e^{z_K} \end{bmatrix}$$

activation function

$$W = \begin{bmatrix} -w_1^T \\ -w_2^T \\ \vdots \\ -w_K^T \end{bmatrix}$$

$$P(Y=0|X, \theta) = \frac{e^{\theta_0^T X}}{e^{\theta_0^T X} + e^{\theta_1^T X}}$$

$$P(Y=1|X, \theta) = \frac{e^{\theta_1^T X}}{e^{\theta_0^T X} + e^{\theta_1^T X}}$$

How do we apply gradient descent?
Optimize the parameters?
We have to use
backpropagation

Cross Entropy Loss function (for N classes) is the softmax loss:

for 2 classes

$$L(W) = - \left[\sum_{i=1}^m \sum_{k=1}^K \mathbb{1}\{y_i=k\} \log \left(\frac{e^{W_k \cdot x_i}}{\sum_{j=1}^K e^{W_j \cdot x_i}} \right) \right]$$

In examples, K classes, $\mathbb{1}\{y_i=k\}$ is an indicator variable which the class prediction made is true. Can be implemented with one-hot encoding with indicator function.

Loss function

cross-entropy and log likelihood loss are basically the same since cross-entropy explicitly computes a softmax calculation and does one-hot encoding, which can be done with an indicator function, and passes it into the log likelihood function.